

TIGHTNESS OF STATIONARY DISTRIBUTIONS OF A FLEXIBLE-SERVER SYSTEM IN THE HALFIN-WHITT ASYMPTOTIC REGIME

BY ALEXANDER L. STOLYAR

Bell Labs, Alcatel-Lucent

We consider a large-scale flexible service system with two large server pools and two types of customers. Servers in pool 1 can only serve type 1 customers, while server in pool 2 are flexible – they can serve both types 1 and 2. (This is a so-called “N-system.” Our results hold for a more general class of systems as well.) The service rate of a customer depends both on its type and the pool where it is served. We study a priority service discipline, where type 2 has priority in pool 2, and type 1 prefers pool 1. We consider the Halfin-Whitt asymptotic regime, where the arrival rate of customers and the number of servers in each pool increase to infinity in proportion to a scaling parameter n , while the overall system capacity exceeds its load by $O(\sqrt{n})$.

For this system we prove tightness of diffusion-scaled stationary distributions. Our approach relies on a single common Lyapunov function $G(x)$, defined on the entire state space as a functional of the *drift-based fluid limits* (DFL). Specifically, $G(x) = \int_0^\infty g(y(t))dt$, where $y(\cdot)$ is the DFL starting at x , and $g(\cdot)$ is a “distance” to the origin. The key part of the analysis is the study of the (first and second) derivatives of the DFLs and function $G(x)$. Our approach, as well as many parts of the analysis, seem quite generic and may be of independent interest.

1. Introduction. In this paper we consider a large-scale service system in the so-called Halfin-Whitt asymptotic regime. Such systems received a lot of attention in the literature, especially in the past 10-15 year, because they find a variety of applications, including, e.g., large customer contact centers [1, 9] and large computer farms in network clouds. The Halfin-Whitt regime, introduced originally in [12], is such that the system capacity (roughly, number of servers) increases in proportion to a scaling parameter n , and exceeds the system load by $O(\sqrt{n})$. It is attractive because it allows – in principle, under a good control algorithm – to achieve both good performance (e.g. waiting times) and high resource utilization.

In the Halfin-Whitt regime, the stochastic process describing the system behavior is usually studied under *diffusion scaling*, i.e. it is centered at the system equilibrium point and scaled down by $n^{-1/2}$. This name reflects the fact that, in the limit on $n \rightarrow \infty$, on any finite time interval, the sequence of diffusion-scaled processes $Y^{(n)}(\cdot)$ “typically” converges to a diffusion process $Y(\cdot)$. Then, a fundamental question is whether or not the following *limit interchange* property holds: *the limit of stationary distributions of $Y^{(n)}(\cdot)$ is equal to the stationary distribution of $Y(\cdot)$* . In turn, the key

AMS 2000 subject classifications: Primary 60K25, 60F17

Keywords and phrases: many server models, drift-based fluid limit, diffusion limit, tightness of stationary distributions, limit interchange, common Lyapunov function

difficulty in establishing the limit interchange property, is verifying the (*stationary distribution*) *tightness* property: *the family of stationary distributions of $Y^{(n)}(\cdot)$ is tight*.

The tightness property in the Halfin-Whitt regime is usually difficult to verify even for systems with single pool of homogeneous servers, if there is more than one type of arriving customers and/or the service time distribution is non-exponential; see [3, 5–7] for the results in this direction. (We note that the problems of verifying the tightness and limit interchange exists not only in the Halfin-Whitt regime, but also in the so-called *conventional heavy traffic* regime; see [2, 8, 10].) More general models, where there are multiple flexible server pools with different capabilities (service rates) w.r.t. different customer types, pose additional challenges. The key additional difficulty is that for such systems the state space is “fractured” into multiple domains, where the process dynamics is very different. Papers [14–16] contain tightness / limit interchange results for some flexible multi-pool models; although, [14, 15] consider a strictly subcritical load regime (different from Halfin-Whitt), in which the capacity exceeds the load by $O(n)$.

One approach for verifying the stationary distribution tightness is to find a single common Lyapunov function, for which an appropriate “negative expected drift” condition can be established. This approach is used in [3, 6, 7, 16]. (Papers [5, 14, 15] use different approaches, not relying on a single Lyapunov function.) Of course, finding/constructing a suitable Lyapunov function is usually the key challenge. For example, paper [3], which proves tightness for a single-pool model with first-come-first-serve discipline and phase-type service time distribution, uses an elaborate *common quadratic Lyapunov function*, of the type proposed in [4]; the tightness result in [3] also requires that the customers waiting in the queue abandon at positive rate. And again, finding single common Lyapunov function is further complicated for flexible multi-pool systems.

1.1. Paper contributions. We consider a flexible multi-pool system with two customer types and two server pools (the so-called “ N -system”), under a priority discipline, in the Halfin-Whitt regime, and prove the stationary distribution tightness result, Theorem 2, which implies the limit-interchange, Corollary 4. (These results hold for a more general class of systems as well, as discussed in Section 6.)

The state space of the diffusion-scaled process for N -system has five domains (where the process drift is given by different affine functions). Nevertheless, we construct a single Lyapunov function $G(x)$ on the entire state space, as a functional of the *drift-based fluid limits* (DFL), which are the deterministic trajectories defined by the drift of the process. Specifically,

$$(1) \quad G(x) = \int_0^\infty g(y(t))dt,$$

where $y(\cdot)$ is the DFL starting at x , and $g(\cdot)$ is a “distance” to the origin. For Lyapunov functions of this type, in a setting more general than needed for the proof of Theorem 2, we give sufficient conditions for the tightness in Theorem 5; the key condition a bound on the Lyapunov function second derivatives. This result may be of independent interest.

The proof of Theorem 2 verifies the conditions of Theorem 5 for the N -system. This requires the analysis of the DFL structure, and of the (first and second) derivatives of DFLs and corresponding functionals $G(x)$ on the initial state x ; it also requires an appropriate choice of the “distance” $g(\cdot)$. Many parts of this analysis are quite generic and may also be of independent interest.

We note that for a *deterministic* dynamic system, with trajectories $y(\cdot)$ defined by a continuous derivative-field, the function $G(x)$ given by (1) is a natural Lyapunov function (as will be illustrated

in Section 2), as long as it is well defined (the integral in (1) is finite). In particular, this observation is used in [13, 17] to establish the existence of a Lyapunov function for stable *deterministic* fluid models. This observation, however, does *not* imply that $G(x)$ defined by (1) via DFLs $y(\cdot)$ can serve as a Lyapunov function for a (family of) *random* processes(es). In this paper we give sufficient conditions under which Lyapunov functions $G(x)$ can be used to establish tightness of stationary distributions, and then verify these conditions for the N -system.

1.2. Layout of the rest of the paper. In Section 2, we informally discuss our general approach and the Lyapunov function construction. Section 3 formally defines the N -system, the Halfin-Whitt regime for it, and states the tightness (Theorem 2) and the limit-interchange (Corollary 4) results. In Section 4, in a setting more general than needed for the N -system, we give a formal construction of the DFLs and the Lyapunov function, and sufficient conditions for the tightness (Theorem 5). Section 5 contains the proof of Theorem 2; here we choose a specific “distance” function g and verify the conditions of Theorem 5 for the N -system. A generalization of the N -system, for which our results still hold, is described in Section 6. Finally, in Section 7, we discuss our approach and results.

1.3. Basic notation. Symbols $\mathbb{R}, \mathbb{R}_+, \mathbb{Z}, \mathbb{Z}_+$ denote the sets of real, real non-negative, integer, and integer non-negative numbers, respectively. In the Euclidean space \mathbb{R}^I (of dimension $I \geq 1$): $|x|$ denotes standard Euclidean norm of vector $x = (x_1, \dots, x_I)$, while $\|x\| = \sum_i |x_i|$ denotes its L_1 -norm; scalar product of two vectors is denoted $x \cdot y = \sum_i x_i y_i$; $\text{diag}(x)$ denotes diagonal square matrix, with diagonal elements given by x ; we write simply 0 for a zero matrix or vector; vectors are written as row-vectors, but in matrix expressions they are viewed as column-vectors (without using a transposition sign). For real numbers u and w : $u \vee w = \max\{u, w\}$, $u \wedge w = \min\{u, w\}$, and $\lfloor u \rfloor$ denotes the largest integer not greater than u .

For a vector-function $y(\cdot) = (y(t), t \geq 0)$, we denote $\|y(\cdot)\| = \sup_{[0, \infty)} \|y(t)\|$. Abbreviation *u.o.c.* means *uniform on compact sets* convergence. If $X(t)$, $t \geq 0$, is a Markov process, we write $X(\infty)$ for a random element with the distribution equal to a stationary distribution of the process. (In all cases considered in this paper, the stationary distribution will be unique.) Symbol \Rightarrow denotes convergence in distribution of random elements; random processes are random elements in the appropriate Skorohod space. For a condition/event H , the indicator function $I\{H\}$ is equal to 1 when H holds and 0 otherwise.

2. The intuition for the Lyapunov function construction. The discussion in this entire section is informal. Consider a deterministic dynamic system governed by ODE

$$(d/dt)y = v(y),$$

where state y is a vector, and the vector-field $v(\cdot)$ is Lipschitz continuous. Suppose the system has unique stable point 0. Let $g(x)$ be a non-negative continuous (and sufficiently smooth) function, which measures a “distance” from 0. (In our results, we will use $g(x)$ which is a smooth approximation of L_1 -norm $\|x\|$.) Suppose that for any initial state $y(0) = x$ the trajectory $y(t)$, $t \geq 0$ converges to 0 and, moreover,

$$(2) \quad G(x) = \int_0^\infty g(y(t))dt < \infty.$$

Then, obviously, $G(\cdot)$ is a Lyapunov function for this dynamic system, in the sense that

$$(d/dt)G(y) = G'(y) \cdot v(y) = -g(y),$$

where G' denotes the gradient of G .

Suppose now that instead of a deterministic system we have a Markov process $Y(\cdot)$, for which vector-field $v(\cdot)$ gives the drift. Then we can define deterministic trajectories $y(\cdot)$, and function $G(\cdot)$, the same way as above. (The trajectories $y(\cdot)$ we call drift-based fluid limits (DFL).) Suppose further that the process generator A is such that

$$(3) \quad AG(y) = G'(y) \cdot v(y) + H(y), \quad |H(y)| \leq C\|G''(y)\|,$$

where C is a constant and G'' denotes the Hessian matrix of second derivatives. (To interpret (3) one can think, for example, of a diffusion process with bounded diffusion coefficients. In this paper we will work *not* with diffusion processes, but rather with diffusion-scaled processes for our queueing system – their behavior can be very different from that of diffusions, especially when the system state is "far" from the equilibrium point. Nevertheless, the process generator will have form (3).) Then, we have

$$AG(y) \leq G'(y) \cdot v(y) + C\|G''(y)\| = -g(y) + C\|G''(y)\|.$$

If we can show that

$$(4) \quad \|G''(y)\| \leq C_1 g(y) + C_2$$

with a sufficiently small C_1 , then for some $\epsilon > 0$,

$$(5) \quad AG(y) \leq -\epsilon g(y) + C_3.$$

This is a Lyapunov-Foster type condition from which we can obtain the steady-state bound $\mathbb{E}g(Y) \leq C_3/\epsilon$.

Finally, suppose we consider a family of processes $Y = Y^{(n)}$, with the drift $v(\cdot)$ and generator A depending on n . If for some common function g such that $g(x) \rightarrow \infty$ as $\|x\| \rightarrow \infty$, we can derive estimates (3)-(5) with constants independent of n , then the family of stationary distributions of $Y = Y^{(n)}$ is tight.

This is the program that we implement in this paper, for the sequence of diffusion-scaled processes for the N -system. The difficult part is obtaining the second derivative bound (4). Since G is defined as a functional of the DFLs $y(\cdot)$, this involves the analysis of the dependence of DFLs on the initial state.

3. N -system with absolute priority. Consider a so-called N -system, with absolute priorities. (See Fig. 1.) There are two customer types, arriving according to as independent Poisson processes with rates $\Lambda_1 > 0$ and $\Lambda_2 > 0$, respectively. There are two server pools, with B_1 and B_2 identical servers, respectively. The total service requirement of any customer is an independent, exponentially distributed random variable with mean 1. A customer of type 2 can only be served by a server in pool 2, and if it does receive service, it does so at rate $\mu_{22} > 0$. A customers of type 1 can be served by a server in either pool 1 or 2, with service rates being $\mu_{11} > 0$ and $\mu_{12} > 0$, respectively. Type 2 customers have absolute (preemptive) priority (in pool 2); namely, if there are X_2 type 2 customers in the system, as many of them as possible, $X_2 \wedge B_2$, receive service in pool 2, and the remaining

$X_2 - X_2 \wedge B_2 = (X_2 - B_2) \vee 0$ wait in the queue. (Here \wedge and \vee denote minimum and maximum, respectively.) Therefore, the total service rate of all type 2 customers is

$$(6) \quad \mu_{22}(X_2 \wedge B_2).$$

The type 1 customers have absolute preference to be served in pool 1, and have lower preempt-resume priority in pool 2. Namely, if there are X_1 type 1 customers in the system, then $X_1 \wedge B_1$ of them are served in pool 1, $[(X_1 - B_1) \vee 0] \wedge [(B_2 - X_2) \vee 0]$ are served in pool 2, and the remaining $[X_1 - (B_1 + B_2) + (X_2 \wedge B_2)] \vee 0$ wait in queue. The total service rate of all type 1 customers is

$$(7) \quad \mu_{11}\{X_1 \wedge B_1\} + \mu_{12}\{[(X_1 - B_1) \vee 0] \wedge [(B_2 - X_2) \vee 0]\}.$$

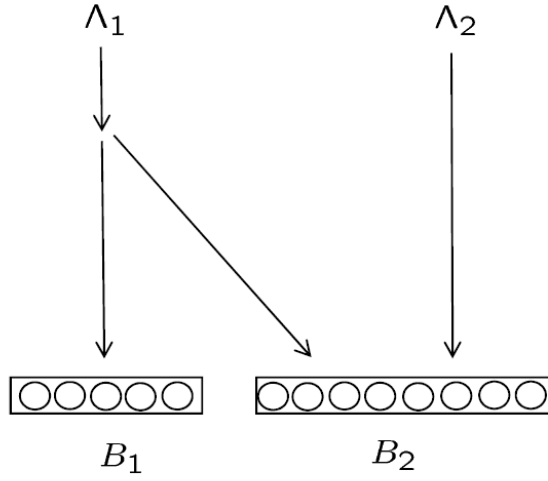


FIG 1. *N-system*.

We consider a sequence of such systems, indexed by a positive scaling parameter n , increasing to infinity. (See Fig. 2.) In a system with parameter n ,

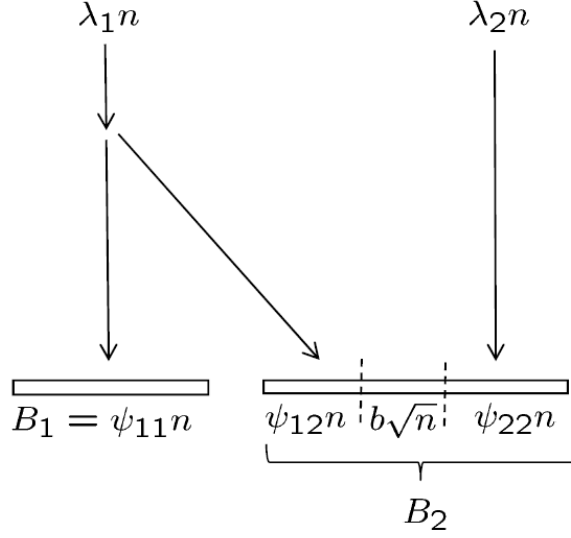
$$(8) \quad \Lambda_1 = \lambda_1 n, \quad \Lambda_2 = \lambda_2 n,$$

$$(9) \quad B_1 = \psi_{11}n, \quad B_2 = \psi_{12}n + \psi_{22}n + b\sqrt{n},$$

where the positive parameters $b, \lambda_1, \lambda_2, \psi_{11}, \psi_{12}, \psi_{22}$ are such that

$$(10) \quad \lambda_2 = \psi_{22}\mu_{22}, \quad \lambda_1 = \psi_{11}\mu_{11} + \psi_{12}\mu_{12}.$$

Given this definition, and the priorities, the system “desired operating point,” which we will refer to as *equilibrium point*, is such that $X_2 = \psi_{22}n$ and $X_1 = \psi_{11}n + \psi_{12}n$, where type 1 customers occupy the entire pool 1 and $\psi_{12}n$ servers in pool 2; the equilibrium point is such that $b\sqrt{n}$ servers in pool 2 are empty – this is the “margin” by which system capacity exceeds its load. (Again, see Fig. 2.)

FIG 2. *N-system in Halfin-Whitt asymptotic regime.*

REMARK 1. To be precise, in the definition of the sequence of systems, we need to make sure that B_1 and B_2 are integer. Equations (9), as written above, assume that B_1 and B_2 "happen to be" integer. We make this assumption throughout the paper to simplify the exposition, while maintaining rigor of the results and arguments. More specifically, we could replace (9) with, for example,

$$(11) \quad B_1 = \lfloor \psi_{11}n \rfloor, \quad B_2 = \lfloor \psi_{12}n + \psi_{22}n + b\sqrt{n} \rfloor.$$

If we do that, it is easy to check that for each n we can choose numbers $\psi_{ij}^{(n)}$, $(ij) = (11), (12), (22)$, and $b^{(n)}$, such that: $|\psi_{ij}^{(n)} - \psi_{ij}| \leq C_{55}/n$ and $|b^{(n)} - b| \leq C_{55}/\sqrt{n}$ for some constant $C_{55} > 0$; (11) can be rewritten as

$$(12) \quad B_1 = \psi_{11}^{(n)}n, \quad B_2 = \psi_{12}^{(n)}n + \psi_{22}^{(n)}n + b^{(n)}\sqrt{n};$$

and (10) can be rewritten as

$$(13) \quad \lambda_2 = \psi_{22}^{(n)}\mu_{22}, \quad \lambda_1 = \psi_{11}^{(n)}\mu_{11} + \psi_{12}^{(n)}\mu_{12}.$$

The sequence of systems will then be defined by (8), (12), (13). Then, the entire analysis in this paper will hold as is, with ψ_{ij} and b replaced everywhere with $\psi_{ij}^{(n)}$ and $b^{(n)}$, respectively. (We note that the components of the equilibrium point, namely $X_2 = \psi_{22}^{(n)}n$ and $X_1 = \psi_{11}^{(n)}n + \psi_{12}^{(n)}n$, need *not* be integer.)

It is easy to see that for each n the process $X^{(n)}(t) = (X_1^{(n)}(t), X_2^{(n)}(t))$, $t \geq 0$, is continuous-time countable irreducible Markov chain, with the state space being (for each n) \mathbb{Z}_+^2 . Further, it is not difficult to check that, for each sufficiently large n , this Markov process is positive recurrent, and therefore has unique stationary distribution. Indeed, due to absolute priority, type 2 customers "do not see" type 1, and therefore $X_2^{(n)}(\cdot)$ in itself is a positive recurrent Markov chain, which in

steady-state occupies on average $\psi_{22}n$ servers in pool 2. This means that on average $\psi_{12}n + b\sqrt{n}$ servers in pool 2 are available to serve type 1 customers; this is in addition to all $\psi_{11}n$ servers on pool 1 which are available exclusively to type 1; therefore, the average total service capacity available to type 1 is

$$\psi_{11}n\mu_{11} + (\psi_{12}n + b\sqrt{n})\mu_{12} = \lambda_1n + b\mu_{12}\sqrt{n} > \lambda_1n.$$

We omit further details, which are rather straightforward.

The diffusion-scaled version $\hat{X}^{(n)}(t) = (\hat{X}_1^{(n)}(t), \hat{X}_2^{(n)}(t))$ of the process $X^{(n)}(t)$ is defined by centering at the equilibrium point and rescaling by $1/\sqrt{n}$:

$$(14) \quad \hat{X}_1^{(n)} = (X_1^{(n)} - \psi_{11}n - \psi_{12}n)/\sqrt{n}, \quad \hat{X}_2^{(n)} = (X_2^{(n)} - \psi_{22}n)/\sqrt{n}.$$

THEOREM 2. *For some $C > 0$ and all sufficiently large n ,*

$$\mathbb{E}\|\hat{X}^{(n)}(\infty)\| \leq C.$$

The proof of Theorem 2 is given in the rest of this paper. It relies on a family of Lyapunov functions (indexed by n), each being a functional of a fluid models, determined by the process drift. Such fluid models will be referred to as a drift-based fluid limits (DFL). In the rest of this section we define DFLs for the N -system under consideration, and give motivation for the form of Lyapunov function. Then, in Section 4, we give the Lyapunov function construction and sufficient tightness conditions (Theorem 5) in a setting that is more general than needed for the N -system. In the following sections we verify the conditions of Theorem 5 for the N -system, thus proving Theorem 2.

For each n , for the unscaled process $X^{(n)}(\cdot)$, we define a drift function (vector field) $V^{(n)} = (V_1^{(n)}, V_2^{(n)})$ for $x = (x_1, x_2) \in \mathbb{R}_+^2$. (Note that it is defined on \mathbb{R}_+^2 , and not just on the lattice \mathbb{Z}_+^2 .) It is defined in the natural way, as the difference of arrival and service rates (see (6)-(7)):

$$(15) \quad V_1^{(n)} = V_1^{(n)}(x) = \Lambda_1 - \mu_{11}\{x_1 \wedge B_1\} - \mu_{12}\{[(x_1 - B_1) \vee 0] \wedge [(B_2 - x_2) \vee 0]\},$$

$$(16) \quad V_2^{(n)} = V_2^{(n)}(x) = \Lambda_2 - \mu_{22}(x_2 \wedge B_2),$$

where $\Lambda_1, \Lambda_2, B_1, B_2$ are the functions of n given in (8)-(10).

Let us denote by L_n the affine mapping $X^{(n)} \rightarrow \hat{X}^{(n)}$, defined by (14). Then, the state space of $\hat{X}^{(n)}$ is $\mathcal{S}^{(n)} \equiv L_n\mathbb{Z}_+^2 \subset \mathcal{X}^{(n)} \equiv L_n\mathbb{R}_+^2 \subset \mathbb{R}^2$. Specifically, $\mathcal{X}^{(n)} = \{x \mid x_1 \geq -\psi_{11}\sqrt{n} - \psi_{12}\sqrt{n}, x_2 \geq -\psi_{22}\sqrt{n}\}$. The drift function for $\hat{X}^{(n)}$ is defined accordingly:

$$v^{(n)}(x) = (1/\sqrt{n})V^{(n)}(L_n^{-1}x), \quad x \in \mathcal{X}^{(n)}.$$

We emphasize, that $v^{(n)}(x)$ is defined on the continuous, convex set $\mathcal{X}^{(n)}$, which contains the discrete state space $\mathcal{S}^{(n)}$. It is important, however, that *at each point $x \in \mathcal{S}^{(n)}$, $v^{(n)}(x)$ gives exactly the average drift of the process.* Namely,

$$(17) \quad v^{(n)}(x) = \sum_{x'} (x' - x)\nu^{(n)}(x, x'),$$

where $\nu^{(n)}(x, x')$ is the Markov process transition rate from state x to state x' ; note that there is only a finite number of "neighbor" states x' for which $\nu^{(n)}(x, x') > 0$.

As $n \rightarrow \infty$, set $\mathcal{X}^{(n)}$ increases and converges to \mathbb{R}^2 .

It is easy to observe that $v^{(n)}(x) = 0$ if and only if $x = 0$; also, uniformly in n , $v^{(n)}(x)$ is Lipschitz continuous. Given Lipschitz continuity of $v^{(n)}$, for any $x \in \mathcal{X}^{(n)}$ there is a unique solution $y^{(n)}(t)$, $t \geq 0$, to the ODE

$$(d/dt)y^{(n)}(t) = v^{(n)}(y^{(n)}(t)), \quad y^{(n)}(0) = x,$$

which is easily seen to stay within $\mathcal{X}^{(n)}$ for all $t \geq 0$. This trajectory $y^{(n)}(t)$, $t \geq 0$, will be called *drift-based fluid limit* (DFL), starting from x .

As we will show later in Section 5.1, each DFL $y^{(n)}(t) \rightarrow 0$ as $t \rightarrow \infty$. Moreover, after a finite time, this convergence is exponentially fast, so that

$$\int_0^\infty \|y^{(n)}(t)\| < \infty.$$

The Lyapunov function we will use to prove Theorem 2 is

$$G^{(n)}(x) = \int_0^\infty g(y^{(n)}(t)) < \infty,$$

where $y^{(n)}(\cdot)$ is the DFL starting from x , and $g(\cdot)$ is a smooth non-negative function (common for all n) approximating $\|\cdot\|$.

REMARK 3. Deterministic trajectories defined by the drift vector field, which we call DFL, have been considered in the literature; see e.g. [3], where they are called fluid models. However, the way we use DFLs in this paper – namely, to directly construct a Lyapunov function from them – is completely different from their use in [3].

3.1. *Limit interchange.* We conclude this section by noting that the tightness of stationary distributions of the processes $\hat{X}^{(n)}(\cdot)$, which follows from Theorem 2, allows us to easily establish the limit interchange result, given in Corollary 4 below. Observe that $v^{(n)}(x) \rightarrow v(x)$ u.o.c. on \mathbb{R}^2 , where $v_2(x) = -\mu_{22}x_2$, and $v_1(x) = -\mu_{12}[x_1 \wedge (b - x_2)]$. (In fact, on any compact set, $v^{(n)}(x) = v(x)$ for all large n .)

COROLLARY 4. *The following convergence holds*

$$(18) \quad \hat{X}^{(n)}(\infty) \Rightarrow \hat{X}(\infty),$$

where $\hat{X}(\cdot)$ is a diffusion process which is a strong solution of SDE

$$(19) \quad d(\hat{X}) = v(\hat{X})dt + (\sigma_1 dW_1, \sigma_2 dW_2),$$

where W_1, W_2 are independent standard Brownian motions and the diffusion coefficients are $\sigma_1 = [\lambda_1 + \psi_{11}\mu_{11} + \psi_{12}\mu_{12}]^{1/2}$, $\sigma_2 = [\lambda_2 + \psi_{22}\mu_{22}]^{1/2}$.

The proof is fairly straightforward, we just give an outline. First, the following convergence on a finite interval holds (see e.g. [11]). Namely, consider a sequence of processes $\hat{X}^{(n)}(\cdot)$ with fixed initial states $\hat{X}^{(n)}(0) \rightarrow x \in \mathbb{R}^2$. Then, for any fixed $T_0 > 0$

$$(20) \quad (\hat{X}^{(n)}(t), t \in [0, T_0]) \Rightarrow (\hat{X}(t), t \in [0, T_0]),$$

where $\hat{X}(\cdot)$ is a strong solution of (19) with initial state $\hat{X}(0) = x$. Then, (18) can be established, together with the existence and uniqueness of a stationary distribution of $\hat{X}(\cdot)$, as follows. We consider the sequence of stationary versions of the processes $\hat{X}^{(n)}(\cdot)$ on a fixed finite time interval $[0, T_0]$, and let $n \rightarrow \infty$. Given tightness of stationary distributions of pre-limit processes, we can choose a subsequence along which $\hat{X}^{(n)}(0) \Rightarrow \tilde{X}(0)$ for some random vector $\tilde{X}(0)$; then we also have $\hat{X}^{(n)}(T_0) \Rightarrow \tilde{X}(0)$. We then use (20) to show that the distribution of $\tilde{X}(0)$ must be a stationary distribution of $\hat{X}(\cdot)$. The uniqueness of the latter stationary distribution is easy to establish, for example, using a coupling argument.

4. Lyapunov function construction and a tightness criterion. The model in this section is quite general (including the N-system as a special case). For this model we define DFLs, construct a functional of DFL, and give sufficient conditions under which this functional can serve as a Lyapunov function to prove tightness of stationary distributions. The section is self-contained, because its main construction and result may be of independent interest. However, it may help the reader to keep the N-system described in Section 3 in mind as an example, to make the material more concrete.

Let $I \geq 1$ be a fixed positive integer. For each $n = 1, 2, 3, \dots$, we consider a Markov chain $\hat{X}^{(n)}(t)$, $t \geq 0$, with a countable state space $\mathcal{S}^{(n)}$ which has the form

$$\mathcal{S}^{(n)} = \{L_n x \mid x \in \mathbb{Z}^I\} \cap \mathcal{X}^{(n)},$$

where $\mathcal{X}^{(n)}$ is a convex closed subset of \mathbb{R}^I , containing 0, and $L_n x = x/\sqrt{n} + s^{(n)}$ with some fixed $s^{(n)} \in \mathbb{R}^I$. The Markov chain is irreducible positive-recurrent and is such that the total transition rate out of any state is upper bounded by $R_1 n$ and any single transition has the jump size of at most R_2/\sqrt{n} , where R_1, R_2 are positive constants independent of n . Defined on $\mathcal{X}^{(n)}$ is a drift function (vector field) $v^{(n)}(x)$, which is Lipschitz continuous uniformly in n . At each point $x \in \mathcal{S}^{(n)}$, $v^{(n)}(x)$ gives exactly the average drift of the process. Namely,

$$(21) \quad v^{(n)}(x) = \sum_{x'} (x' - x) \nu^{(n)}(x, x'),$$

where $\nu^{(n)}(x, x')$ is the Markov process transition rate from state x to state x' ; note that there is only a finite number of "neighbor" states x' for which $\nu^{(n)}(x, x') > 0$.

For any $x \in \mathcal{X}^{(n)}$, there is a unique solution $y^{(n)}(t)$, $t \geq 0$, to the ODE

$$(d/dt)y^{(n)}(t) = v^{(n)}(y^{(n)}(t)), \quad y^{(n)}(0) = x,$$

which stays within $\mathcal{X}^{(n)}$. This solution is called *drift-based fluid limit* (DFL), starting from x .

Suppose a continuous non-negative function $g(x)$, $x \in \mathbb{R}^I$, is fixed. For $x \in \mathcal{X}^{(n)}$ define

$$(22) \quad G^{(n)}(x) = \int_0^\infty g(y^{(n)}(t)) dt, \quad y^{(n)}(0) = x,$$

where $y^{(n)}(\cdot)$ is the DFL starting from x .

Denote by $\nabla_z G^{(n)}(x)$ the directional derivative of $G^{(n)}$ at $x \in \mathcal{X}^{(n)}$ in the direction of vector $z \in \mathbb{R}^I$:

$$\nabla_z G^{(n)}(x) \doteq \lim_{\delta \downarrow 0} \frac{1}{\delta} [G^{(n)}(x + z\delta) - G^{(n)}(x)],$$

when the limit exists. (To be precise, if x is on the boundary of $\mathcal{X}^{(n)}$, it is also required that the direction z from x points into $\mathcal{X}^{(n)}$.) Then, $\nabla_{z_*}[\nabla_z G^{(n)}](x)$ is the second derivative, first in the direction z and then z_* .

THEOREM 5. *Suppose that for any $C_1 > 0$, there exists a function $g(x)$, $x \in \mathbb{R}^I$, and a constant $C_2 > 0$, such that the following conditions hold uniformly in n .*

- (i) *$g(x)$ is Lipschitz continuous non-negative and such that $g(x) \rightarrow \infty$ as $x \rightarrow \infty$.*
- (ii) *Function $G^{(n)}(x)$, $x \in \mathcal{X}^{(n)}$, is finite for all x ; it has continuous gradient $\nabla G^{(n)}(x)$; for any x and any fixed unit-length vectors $z, z_* \in \mathbb{R}^I$,*

$$(23) \quad \limsup_{\delta \downarrow 0} \frac{1}{\delta} \left| \nabla_z G^{(n)}(x + z_* \delta) - \nabla_z G^{(n)}(x) \right| \leq C_1 g(x) + C_2;$$

$$(24) \quad G^{(n)}(x) \rightarrow \infty, \quad x \rightarrow \infty.$$

Then, for some $C_3 > 0$, uniformly in n ,

$$\mathbb{E}g(\hat{X}^{(n)}(\infty)) \leq C_3.$$

The second derivative condition (23) is the key one. It implies that $|\nabla_{z_*}[\nabla_z G^{(n)}](x)| \leq C_1 g(x) + C_2$ if this second derivative exists. An equivalent form of (23) is as follows: for any compact set $D \subseteq \mathcal{X}^{(n)}$ and any unit-length vector $z \in \mathbb{R}^I$, the first derivative $\nabla_z G^{(n)}$ within D is Lipschitz continuous with constant

$$C_1 \max_D g(x) + C_2.$$

Proof of Theorem 5. By definition of $G^{(n)}$ and its assumed continuous differentiability,

$$(25) \quad \nabla G^{(n)}(x) \cdot v^{(n)}(x) = \nabla_{v^{(n)}(x)} G^{(n)}(x) = -g(x).$$

Let $A^{(n)}$ denote the infinitesimal generator of the Markov process $X^{(n)}$. For any fixed $k > 0$, the function $G^{(n),k} \doteq G^{(n)} \wedge k$ is such that it has constant value k for all states $x \in \mathcal{S}^{(n)}$ except a finite subset. Then it is easy to see that function $G^{(n),k}$ is within the domain of $A^{(n)}$ and

$$(26) \quad \mathbb{E}A^{(n)}G^{(n),k}(\hat{X}^{(n)}(\infty)) = 0.$$

(See also [7], page 31, for more details.) For all $x \in \mathcal{S}^{(n)}$ we have

$$A^{(n)}G^{(n),k}(x) \leq \nabla G^{(n)}(x) \cdot v^{(n)}(x) + r^{(n)}(x)(1/2)h^{(n)}(x)(R_2/\sqrt{n})^2,$$

where R_2/\sqrt{n} is the maximum possible size of one jump of the process, $r^{(n)}(x) \leq R_1 n$ is the total transition rate from state x , and the second-term coefficient $h^{(n)}(x)$ is bounded as $|h^{(n)}(x)| \leq C_1 g(x) + 2C_2 < \infty$. Recalling also (25), we obtain

$$A^{(n)}G^{(n),k}(x) \leq -g(x) + (1/2)R_1 R_2^2 [C_1 g(x) + 2C_2].$$

We choose a sufficiently small constant $C_1 > 0$ (and then a corresponding function g and constant C_2), so that

$$A^{(n)}G^{(n),k}(x) \leq -\epsilon g(x) + C'_2, \quad \text{for some } \epsilon > 0, \quad C'_2 > 0.$$

Obviously, if $x \in \mathcal{S}^{(n)}$ is such that $G^{(n),k}(x) = k$, which is equivalent to $G^{(n)}(x) \geq k$, then

$$A^{(n)}G^{(n),k}(x) \leq 0.$$

Therefore, from (26) we obtain

$$\mathbb{E}[-\epsilon g(\hat{X}^{(n)}(\infty)) + C'_2]I\{G^{(n)}(\hat{X}^{(n)}(\infty)) < k\} \geq 0,$$

or

$$\mathbb{E}g(\hat{X}^{(n)}(\infty))I\{G^{(n)}(\hat{X}^{(n)}(\infty)) < k\} \leq C'_2/\epsilon.$$

Letting $k \rightarrow \infty$, by monotone convergence,

$$\mathbb{E}g(\hat{X}^{(n)}(\infty)) \leq C'_2/\epsilon.$$

The constant in the RHS is same for all n . \square

5. Proof of Theorem 2. We will prove Theorem 2 by choosing specific function $g(\cdot)$ and then verifying (in Theorem 10) the assumptions of Theorem 5 for N -system.

In this section, we study properties of DFL trajectories and their $G^{(n)}$ -functionals, for a system with a fixed scaling parameter n . We will drop upper index (n) from now on. So, for example, will write simply \mathcal{X} and $y(t)$ instead of $\mathcal{X}^{(n)}$ and $y^{(n)}(t)$, respectively. (However, the expressions may contain n as a variable.) From this point on in the paper, we say that C is a *universal constant* if C depends only on the system parameters $\lambda_i, \psi_{ij}, \mu_{ij}, b$, but does *not* depend on scaling parameter n . (If the sequence of systems is defined as in Remark 1, then a universal constant C depends on the values ψ_{ij} and b , and *not* on the sequences $\psi_{ij}^{(n)}$ and $b^{(n)}$.)

5.1. Basic DFL properties. First derivatives of DFLs and the Lyapunov function. In this subsection we first establish some basic properties of DFLs and their directional (Gateaux) derivatives. Then we specify function $g(\cdot)$, and obtain the expressions for the first derivatives of the corresponding function $G(\cdot)$. (All results of this subsection hold for systems far more general than N -system. In particular, they still hold for the systems under the *Leaf Activity Priority* LAP discipline in [14, 15], in the Halfin-Whitt regime; our priority discipline for the N -system is a special case of LAP.)

The DFL trajectories $y(\cdot)$ have the following structure. Recall that $v(x)$ is (uniformly in n) Lipschitz continuous on the entire \mathcal{X} . There is a finite number M (same for any n) of domains, indexed by $m = 0, \dots, M-1$; within each of them $v(x)$ is a given linear function. More precisely, the DFL satisfies a linear ODE

$$(d/dt)y(t) = v(y(t)) = u^m y(t) + a^m,$$

where u^m is a constant $I \times I$ matrix (same for each n), and a^m is a constant vector (depending on n). Informally speaking, a domain is determined by which service pools a fully occupied and/or which queues are non-empty.

Formally, the domains are easier to define (and think of) in terms of unscaled quantities $X_1 \geq 0$ and $X_2 \geq 0$, and unscaled pool sizes $B_1 = \psi_{11}n$ and $B_2 = \psi_{12}n + \psi_{22}n + b\sqrt{n}$. Each domain is defined by a combination of the directions of three strict inequalities:

$$(27) \quad X_1 < B_1 \quad \text{or} \quad X_1 > B_1,$$

$$(28) \quad X_2 < B_2 \text{ or } X_2 > B_2,$$

$$(29) \quad X_1 + X_2 < B_1 + B_2 \text{ or } X_1 + X_2 > B_1 + B_2.$$

However, we exclude two combinations, or conditions, $(X_1 < B_1, X_2 < B_2, X_1 + X_2 > B_1 + B_2)$ and $(X_1 > B_1, X_2 > B_2, X_1 + X_2 < B_1 + B_2)$, because they produce the empty set; and we replace ("merge") the conditions $(X_1 < B_1, X_2 > B_2, X_1 + X_2 > B_1 + B_2)$ and $(X_1 < B_1, X_2 > B_2, X_1 + X_2 < B_1 + B_2)$ into one condition $(X_1 < B_1, X_2 > B_2)$ because this condition alone determines the form of $v(x)$. So, there are $M = 5$ domains in total. The diffusion-scaling mapping L_n , defined by (14), transforms them into 5 (diffusion-scale) domains, denoted $\mathcal{X}^0, \dots, \mathcal{X}^4$. Note that the domains are defined by strict inequalities, so they do not cover the entire space \mathcal{X} . The domain closures are $\bar{\mathcal{X}}^1, \dots, \bar{\mathcal{X}}^5$, these do cover the entire \mathcal{X} . By these definitions, if a point belongs to the intersection of the closures of more than one domain, then necessarily at least one of the equalities (in terms of unscaled quantities), $X_1 = B_1, X_2 = B_2, X_1 + X_2 = B_1 + B_2$, holds.

In particular, consider the unscaled domain $(X_1 > B_1, X_2 < B_2, X_1 + X_2 < B_1 + B_2) = (X_1 > B_1, X_1 + X_2 < B_1 + B_2)$; it is such that there are no queues and pool 1 fully occupied. The corresponding diffusion-scaled domain is $\mathcal{X}^0 = \{x \in \mathcal{X} \mid x_1 > -\psi_{12}\sqrt{n}, x_1 + x_2 < b\}$. In this domain $v(x) = (-\mu_{12}x_1, -\mu_{22}x_2)$, i.e. $u^0 = \text{diag}(-\mu_{12}, -\mu_{22})$ and $a^0 = 0$, and therefore the components y_1 and y_2 evolve independently. Moreover, there exists a universal constant $\alpha > 0$, such that if $y(t)$ starts from a point $y(0) \in \mathcal{X}^{0,\alpha} \doteq \{\|x\| \leq \alpha\} \subset \mathcal{X}^0$, then $y(t)$ never leaves domain \mathcal{X}_0 , which in turn means that the trajectory is simply

$$y_i(t) = y_i(0)e^{-\mu_{i2}t}, \quad i = 1, 2.$$

From now such constant α and the corresponding sub-domain $\mathcal{X}^{0,\alpha}$ will be fixed.

Consider one more unscaled domain $(X_1 > B_1, X_2 < B_2, X_1 + X_2 > B_1 + B_2)$. Here, pool 1 is fully occupied by type 1, pool 2 is fully occupied by X_2 type 2 customers and $B_2 - X_2$ type 1 customers, and $X_1 - B_1 - (B_2 - X_2) = X_1 + X_2 - B_1 - B_2 > 0$ type 1 customers waiting in the queue. On the diffusion scale, the domain (let us label it $m = 1$) is: $\mathcal{X}^1 = \{x \in \mathcal{X} \mid x_1 > -\psi_{12}\sqrt{n}, x_2 < \psi_{12}\sqrt{n} + b, x_1 + x_2 > b\}$, and we have

$$v(x) = ((-b + x_2)\mu_{12}, -\mu_{22}x_2),$$

with the corresponding u^1 and a^1 . For the remaining 3 domains the $v(x)$ is determined similarly.

The equations for a DFL $y(\cdot)$ can be summarized as follows. The trajectory of y_2 is not affected by y_1 and satisfies ODE

$$(30) \quad (d/dt)y_2 = -\mu_{22}[y_2 \wedge (\psi_{12}\sqrt{n} + b)].$$

If $y_1 \leq -\psi_{12}\sqrt{n}$ (which corresponds to unscaled condition $X_1 \leq B_1$),

$$(31) \quad (d/dt)y_1 = -\mu_{11}(y_1 + \psi_{12}\sqrt{n}) + \mu_{12}\psi_{12}\sqrt{n} \geq \mu_{12}\psi_{12}\sqrt{n}.$$

If $y_1 \geq -\psi_{12}\sqrt{n}$ ($X_1 \geq B_1$) and $y_2 \geq \psi_{12}\sqrt{n} + b$ ($X_2 \geq B_2$),

$$(32) \quad (d/dt)y_1 = \mu_{12}\psi_{12}\sqrt{n}.$$

If $y_1 \geq -\psi_{12}\sqrt{n}$ ($X_1 \geq B_1$), $y_2 \leq \psi_{12}\sqrt{n} + b$ ($X_2 \leq B_2$), and $y_1 + y_2 \leq b$ ($X_1 + X_2 \leq B_1 + B_2$), that is in domain $\bar{\mathcal{X}}^0$,

$$(33) \quad (d/dt)y_1 = -y_1\mu_{12}.$$

If $y_1 \geq -\psi_{12}\sqrt{n}$ ($X_1 \geq B_1$), $y_2 \leq \psi_{12}\sqrt{n} + b$ ($X_2 \leq B_2$), and $y_1 + y_2 \geq b$ ($X_1 + X_2 \geq B_1 + B_2$), that is in domain $\bar{\mathcal{X}}^1$,

$$(34) \quad (d/dt)y_1 = (-b + y_2)\mu_{12}.$$

For a given fluid trajectory, let us call time point $t \geq 0$ a *switching point* if $y(t)$ belongs to the intersection of two or more closed domains $\bar{\mathcal{X}}^m$. (i.e. it is on a boundary separating different domains).

LEMMA 6. *For some universal constants $T > 0$, $C' > 0$ and (integer) $K' > 0$, DFL trajectories $y(\cdot)$ satisfy the following conditions. It is always assumed that $y(0) \in \mathcal{X}$, and we denote $y(0) = x$.*

(i) *Let $\tau \geq 0$ be the first time a DFL reaches set $\mathcal{X}^{0,\alpha}$. Then, $\tau \leq T\|x\|$. In addition, $\|y(\cdot)\| \leq C'\|x\|$.*

(ii) *DFL $y(\cdot)$ depends on its initial state x continuously, in the sense of $\|y(\cdot)\|$ -norm.*

(iii) *DFL $y(\cdot)$ has at most K' switching points, $t_1 < t_2 < \dots < t_K$, $0 \leq K \leq K'$, and $t_K < \|x\|T$. Moreover, the set of switching points is upper semicontinuous: as a DFL initial state converges to x , the limiting points of the set of switching points are within the set of switching points for initial state x .*

(iv) *For any interval $[C_3, C_4]$, not containing 0, there exists a constant $T_3 > 0$ (independent of n), such that the total time the condition $y_i(t) \in [C_3, C_4]$ holds for at least one i , is upper bounded by T_3 .*

Proof of Lemma 6. Given equation (30), condition $y_2(t) = \psi_{12}\sqrt{n} + b$ ($X_2 = B_2$) can hold at most at one point $t_2 \geq 0$, which will be a switching point. Similarly, by (31), there is at most one point $t_1 \geq 0$, at which condition $y_1(t) = -\psi_{12}\sqrt{n}$ (corresponding to $X_1 = B_1$) can hold, and if so, it will be a switching point.

Denote $t' = t_1 \vee t_2$. It is easy to see that for some universal constant $C_{30} > 0$,

$$(35) \quad t' \leq C_{30}\|x\|, \quad \|y(t')\| \leq C_{30}\|x\|.$$

Indeed, $|y_2(t)|$ is non-increasing in $[0, \infty)$, and $t_2 \leq |x_2|/[(\psi_{12}\sqrt{n} + b)\mu_{22}] \leq |x_2|/[\psi_{12}\mu_{22}\sqrt{n}]$. In the interval $[0, t_1]$, $y_1(t)$ is negative non-decreasing, and then $|y_1(t)|$ is non-increasing; and $t_1 \leq |x_1|/[\psi_{12}\mu_{12}\sqrt{n}]$. If $t_2 > t_1$, then in the interval $[t_1, t_2]$, $(d/dt)y_1(t) = \psi_{12}\mu_{12}\sqrt{n}$, and therefore $|y_1(t_2) - y_1(t_1)| \leq \psi_{12}\mu_{12}\sqrt{n}t_2$; given the bound on t_2 , we see that $|y_1(t_2) - y_1(t_1)|$ is upper bounded by $|x_2|$ times a universal constant. These observations imply (35).

For all $t > t'$, conditions $y_2(t) < \psi_{12}\sqrt{n} + b$ ($X_2 < B_2$) and $y_1(t) > -\psi_{12}\sqrt{n}$ ($X_1 > B_1$) hold. Therefore, $y(t)$ can be only in one of the two domains $\bar{\mathcal{X}}^0$ or $\bar{\mathcal{X}}^1$, depending on whether $y_1 + y_2 \leq b$ (no queues) or $y_1 + y_2 \geq b$ (queue size $y_1 + y_2 - b$ of type 1). It is easy to see from equations $(d/dt)y_2 = -\mu_{22}y_2$, (33), (34), that if $y(t)$ is in $\bar{\mathcal{X}}^1$, then the trajectory eventually leaves $\bar{\mathcal{X}}^1$ and can never return. This implies that at most two transitions between \mathcal{X}^0 and \mathcal{X}^1 can occur after t' . Specifically, either the trajectory stays in \mathcal{X}^0 , or it is in \mathcal{X}^1 and then \mathcal{X}^0 , or it is in \mathcal{X}^0 then \mathcal{X}^1 then \mathcal{X}^0 . The boundary cases are also possible; for example, the trajectory may stay in the open domain

\mathcal{X}^0 at all times, except at exactly one point $t \geq t'$ it "touches" the boundary, i.e. $y_1 + y_2 = b$. To summarize, after t' there are at most two switching points.

Denote by t'' the first time $t \geq t'$ when $\|y_2(t)\| \leq \alpha/4$. We have $t'' - t' = 0 \vee (1/\mu_{22}) \log[\|y_2(t')\|/(\alpha/4)] \leq C_{31}\|x\| + C_{32}$, for some universal C_{31} and C_{32} . (C_{32} depends on α , which in turn is universal.) In the interval $[t', t'']$ the value of $|y_1|$ cannot increase by more than $C_{33}|y_2(t')| \leq C_{34}\|x\|$, for universal $C_{33}, C_{34} > 0$. (If $y_1 \leq 0$, then $(d/dt)y_1 \geq 0$. If $y_1 \geq 0$, then $(d/dt)y_1 \leq \mu_{12}|y_2|$, and recall that $(d/dt)y_2 = -\mu_{22}y_2$.) Therefore, $|y_1(t'')| \leq C_{35}\|x\|$, for universal $C_{35} > 0$. Starting t'' , if type 1 has non-zero queue, $(d/dt)|y_1| = (d/dt)y_1 \leq -C_{36} < 0$, for universal $C_{36} > 0$; and if type 1 does not have queue, then $(d/dt)|y_1| = -\mu_{12}|y_1|$. Consider the first time $t''' \geq t''$ when $|y_1| \leq \alpha/4$. We conclude that $t''' \leq T\|x\| + C_{37}$ and $\sup_{[0, t''']} \|y(t)\| \leq C'\|x\|$ for some universal positive constants T, C', C_{37} . Obviously, $t''' \geq \tau$, so that $\tau \leq T\|x\| + C_{37}$. However, if $\|x\| \leq \alpha$, i.e. $y(0) = x$ is already in $\mathcal{X}^{0,\alpha}$, then obviously $\tau = 0$. Therefore, in the bound $\tau \leq T\|x\| + C_{37}$, we can drop C_{37} by rechoosing T , if necessary.

For future reference, we also make the following observation. *Suppose, $\mu_{12} = \mu_{22}$. Then, there can be at most one switching point after time t' , let us call it $t_3 \geq t'$, and it is such that $y(t) \in \mathcal{X}^0$ for all $t > t_3$. Indeed, in this case, in the domain $\bar{\mathcal{X}}^0$, we have simply $(d/dt)[y_1 + y_2] = -\mu_{22}[y_1 + y_2]$.*

Let us prove properties (i)-(iv). In fact, (i) has been proved already. For a given x , let us choose τ' such that $\tau < \tau'$ for all initial states sufficiently close to x . (On a finite interval $[0, \tau']$, $y(\cdot)$ depends on the initial state continuously, because it is a solution to an ODE with Lipschitz continuous RHS.) But, for $t \geq \tau'$, the DFL with any initial state close to x is such that $y(t) \in \mathcal{X}^{0,\alpha}$; this implies uniform convergence across all $t \geq 0$, which proves (ii). The part of property (iii), stating that there is at most K' switching points, all of which are smaller than $\tau \leq T\|x\|$, has already been proved, in fact we specified that $K' \leq 4$. Then, the upper continuity of the set of switching points follows from continuity of trajectories w.r.t. initial state; this proves (iii). Consider a fixed interval $[C_3, C_4]$, not containing 0. It is clear from (30) that $y_2(t)$ can spend only a finite time within $[C_3, C_4]$. Now, $y_1(t)$ can be in $[C_3, C_4]$ only after time t_1 , and then in every domain the trajectory visits $y_1(t)$ satisfies one of the equations (32)-(34). If we examine each of these equations (and recall that (34) holds within domain $\bar{\mathcal{X}}^1$, where $(d/dt)y_2 = -\mu_{22}y_2$), we see that *even if the equation were to hold up to infinite time*, $y_1(t)$ can spend only a finite time within $[C_3, C_4]$. And there is only a finite, uniformly bounded number of domains that a trajectory can visit. This proves (iv). \square

Next, let us consider the first-order dependence of DFL on the initial state. Let $y(t; x)$ denote $y(t)$ with initial state $y(0) = x \in \mathcal{X}$. For any $x \in \mathcal{X}$ and any direction $z \in \mathbb{R}^I$ (which does not point outside \mathcal{X}), we use the following notation for the directional (Gateaux) derivative of $y(t; x)$ at x in the direction z :

$$\nabla_z y(\cdot; x) \doteq \lim_{\delta \downarrow 0} \frac{1}{\delta} [y(\cdot; x + z\delta) - y(\cdot; x)].$$

THEOREM 7. (i) *For any fixed $x \in \mathcal{X}$ and a fixed vector z , the directional derivative*

$$\xi(\cdot) = \xi(\cdot; x, z) = \nabla_z y(\cdot; x)$$

exists. It has the following structure. Let $0 < t_1 < t_2 < \dots < t_K$ be the switching points of $y(\cdot; x)$. Then, $\xi(0) = z$, and in each interval $[0, t_1], [t_1, t_2], \dots, [t_K, \infty)$, ξ satisfies linear homogeneous ODE

$$(d/dt)\xi = u^m \xi,$$

where matrix u^m is the matrix u for the domain $\bar{\mathcal{X}}^m$ containing $y(t; x)$.
Solutions $q(t)$, $t \geq 0$, to the equation $(d/dt)q = u^m q$, for any m , are such that

$$(36) \quad \|q(\cdot)\| \leq C_7 \|q(0)\|$$

for a universal constant $C_7 > 0$.

(ii) The derivative $\xi(\cdot; x, z)$ depends on (x, z) continuously.

(iii) There exists a universal constant $C_8 > 0$, such that

$$\|\xi(\cdot; x, z)\| \leq C_8 \|\xi(0; x, z)\| = C_8 \|z\|.$$

Proof. The proof of (i) relies on the following observations.

(a) In any time interval, where both $y(t; x + z\delta)$ and $y(t; x)$ are within same domain $\bar{\mathcal{X}}^m$, they are governed by the same ODE $(d/dt)y = v^m(y)$, and therefore their difference $\Delta y(t) = y(t; x + z\delta) - y(t; x)$, is governed by the linear homogeneous ODE $(d/dt)\Delta y = u^m \Delta y$. Moreover, it is easy to check that within any domain $\bar{\mathcal{X}}^m$ the corresponding matrix u^m is such that $\|\Delta y(t)\|$ can increase at most by some universal factor C_8 . Indeed, consider Δy_2 first, and then Δy_1 . The equation for Δy_2 is either

$$(37) \quad (d/dt)\Delta y_2 = -\mu_{22}\Delta y_2$$

or

$$(38) \quad (d/dt)\Delta y_2 = 0;$$

in either case $|\Delta y_2|$ can increase at most by a factor (in fact, it cannot increase). The equation for Δy_1 is

$$(d/dt)\Delta y_1 = u_{11}^m \Delta y_1 + u_{12}^m \Delta y_2,$$

where $u_{11}^m = 0$ or $u_{11}^m = -\mu_{11}$ or $u_{11}^m = -\mu_{12}$; we also note that if Δy_2 satisfies (38) then necessarily $u_{12}^m = 0$. We see that in any case, in any time interval, $|\Delta y_2(t)|$ is upper bounded by the initial $\|\Delta y_2\|$ times a universal constant. This observation, in particular, proves (36).

(b) The total length of "switching intervals", where $y(t; x + z\delta)$ and $y(t; x)$ belong to different domains vanishes as $\delta \rightarrow 0$ (by upper semicontinuity of the set of switching points), and therefore the total change of $\Delta y(t)$ within those intervals is "small". More precisely, let t be fixed and $[\theta_1, \theta_2]$ be a switching interval such that $\theta_1, \theta_2 \rightarrow t$. Then, $\|\Delta y(\theta_2) - \Delta y(\theta_1)\| / \|\Delta y(\theta_1)\| \rightarrow 0$, because $v(x)$ is Lipschitz.

Combining observations (a) and (b), and further observing that the number of intervals where both $y(t; x + z\delta)$ and $y(t; x)$ are within same domain $\bar{\mathcal{X}}^m$ (i.e. outside the switching intervals) is upper bounded, we take the $\delta \downarrow 0$ limit to obtain (i).

(ii) This follows from the continuity of the set of switching points on x .

(iii) By (36), in any domain $\|\xi(t)\|$ can increase at most by some factor C_7 . There is only a finite number of domains that $y(t)$ visits. This proves (iii). \square

We now introduce a specific function g , which we will use in the definition (22) of the Lyapunov function W .

DEFINITION 8. Let parameter $C > 0$ be fixed. Let a function $f(\eta)$ of real η be fixed, which satisfies the following conditions. It is a non-negative, even, convex, twice continuously differentiable, $f(\eta) = 0$ for $\eta \in [-C, C]$, $f'(\eta) = -1$ for $\eta \leq -C - 1$, $f'(\eta) = 1$ for $\eta \geq C + 1$. (Such a function can be defined explicitly. Since C is a parameter, essentially, we fix the shape of function $f(C + \zeta)$, $\zeta \geq 0$.) Note that both f' and f'' are uniformly bounded, and $f'' = 0$ outside of the intervals $[-C - 1, -C]$ and $[C, C + 1]$. Then, let

$$g(x) = \sum_i f(x_i).$$

Obviously, $|f(\eta) - |\eta||$ is uniformly bounded by a constant, and then so is $|g(x) - \|x\||$.

Then, by (22) we have $G(x) = \sum_i G_i(x)$, where

$$(39) \quad G_i(x) = \int_0^\infty f(y_i(t)) dt, \quad y(0) = x.$$

THEOREM 9. For $i = 1, 2$, the following holds. For any $x \in \mathcal{X}$ and any direction vector z ,

$$(40) \quad \nabla_z G_i(x) = \int_0^\infty f'(y_i(t; x)) \xi_i(t; x, z) dt.$$

Function $\nabla_z G_i(x)$ is continuous in (x, z) .

Proof. Expression (40) follows from Theorem 7(i) and the fact that f' is continuous bounded. The continuity of $\nabla_z G_i(x)$ is obtained using Theorem 7(ii). \square

5.2. Second derivative bounds for the Lyapunov function.

THEOREM 10. The assumptions of Theorem 5 hold for a function g in Definition 8, with appropriately chosen parameter $C > 0$.

Note that for a function g satisfying Definition 8, condition (i) of Theorem 5 holds automatically. Condition (24) is also automatic given the definition of G and basic properties of DFL, namely the fact that the time for a DFL to reach a given compact set increases to infinity as $x \rightarrow \infty$. Therefore, to prove Theorem 10, it remains to prove condition (23), and it suffices to prove it separately for G_i , $i = 1, 2$ (see (39)). We will do this first for the case $\mu_{22} \neq \mu_{12}$, and then for $\mu_{22} = \mu_{12}$. (The proof of condition (23) in this section applies to the N-system, as well as its generalization described in Section 6. It does *not* apply for LAP discipline.)

For a given x and a time $\tau^* > 0$, denote by $S(\tau^*; x)$ the set of time points, consisting of τ^* and all switching points $0 \leq t < \tau^*$ of the DFL $y(\cdot; x)$.

LEMMA 11. Suppose $\mu_{22} \neq \mu_{12}$. For any $\epsilon > 0$ there exists a sufficiently large $C_9 > 0$, such that, for all sufficiently large n , the following holds for any fixed x and any unit-length vector z . Let τ_9 be the first time the DFL $y(\cdot; x)$ hits set $\{\|y\| \leq C_9\}$. Then for all sufficiently small $\delta > 0$, any point in $S(\tau_9; x + z\delta)$ is within distance at most $\epsilon\delta$ from a point in $S(\tau_9; x)$.

Proof. Consider a switching point $t \in S(\tau_9; x)$ of DFL $y(\cdot) = y(\cdot; x)$. By definition of τ_9 , it is such that $\|y(t)\| \geq C_9$. The switching point is on the boundary of multiple domain closures, and therefore one or more equalities

$$(41) \quad y_1(t) = -\psi_{12}\sqrt{n}, \quad y_2(t) = \psi_{12}\sqrt{n} + b, \quad y_1(t) + y_2(t) = b,$$

defining the domain boundaries, hold. If the first or second equality holds, then $|y'_i(t)|$ is large for large n . If $y_1(t) + y_2(t) = b$, then for t to be a switching point, it is necessary that $y(t) \in \bar{X}^0$; then $y'_1(t) + y'_2(t) = -\mu_{12}y_1(t) - \mu_{22}y_2(t) = -(\mu_{12} - \mu_{22})y_1(t) - \mu_{22}b$; conditions $y_1(t) + y_2(t) = b$ and $\|y(t)\| = |y_1(t)| + |y_2(t)| \geq C_9$ imply that if C_9 is large then so is $|y_1(t)|$, and then $|y'_1(t) + y'_2(t)|$ is large as well. We conclude that if any of the three equalities (41) holds, then for all $n \geq n'$ we have $|y'_1(t)| \geq C_{19}$ or $|y'_2(t)| \geq C_{19}$ or $|y'_1(t) + y'_2(t)| \geq C_{19}$, respectively, where the constant $C_{19} > 0$ can be made arbitrarily large by choosing large enough n' and C_9 . This means that, first, the domains in which the trajectory $y(\cdot; x)$ is in before and after the switching point t are uniquely defined. Second, since the distance between $y(\cdot; x + z\delta)$ and $y(\cdot; x)$ does not exceed $C_{20}\delta$ at all times, where $C_{20} > 0$ is a universal constant (this follows from Theorem 7), and $v(\cdot)$ is Lipschitz, any point in $S(\tau_9; x + z\delta)$ must be within $2C_{20}\delta/C_{19}$ of a point in $S(\tau_9; x)$. Since C_{20} is universal and C_{19} can be made arbitrarily large (by choosing C_9 large), the result follows. \square

Recall that to prove Theorem 10, it remains to prove the second derivative condition (23). The proper second derivative may not exist, hence we must “settle” for the estimate (23). But, to illustrate the proof that follows, let us write down the expression for the second derivative, by formally applying ∇_{z*} differentiation to (40) (this expression is *not* used in the proof):

$$(42) \quad \nabla_{z*} \nabla_z G_i(x) = \int_0^\infty f''(y_i(t; x)) \xi_i(t; x, z_*) \xi_i(t; x, z) dt +$$

$$(43) \quad \int_0^\infty f'(y_i(t; x)) \nabla_{z*} \xi_i(t; x, z) dt.$$

Proof of Theorem 10, case $\mu_{22} \neq \mu_{12}$. We choose small $\epsilon > 0$ and then $C_9 > 0$ as in Lemma 11. Then choose parameter $C > 0$ of function g large enough so that any DFL starting from the set $\{\|y\| \leq 2C_9\}$ never hits set $\{\|y\| \geq C\}$. (We can do this by Lemma 6(i).)

For $i = 1, 2$ consider

$$(44) \quad \begin{aligned} & \frac{1}{\delta} [f'(y_i(t; x + z_*\delta)) \xi_i(t; x + z_*\delta, z) - f'(y_i(t; x)) \xi_i(t; x, z)] \\ &= \frac{1}{\delta} [f'(y_i(t; x + z_*\delta)) - f'(y_i(t; x))] \xi_i(t; x, z) + \end{aligned}$$

$$(45) \quad \frac{1}{\delta} f'(y_i(t; x + z_*\delta)) [\xi_i(t; x + z_*\delta, z) - \xi_i(t; x, z)].$$

(The integrals of the terms (44) and (45), correspond to the integrals (42) and (43), respectively, in the formal second derivative expression.)

Since $f(\cdot)$ has bounded second derivative, the term (44) converges (uniformly in t) to

$$f''(y_i(t; x))\xi_i(t; x, z_*)\xi_i(t; x, z).$$

The integral of this over $t \in [0, \infty)$ is bounded because the total time any trajectory spends in the set $\{C \leq \|y_i\| \leq C + 1\}$ is uniformly bounded (by Lemma 6(iv).)

In the term (45), $f'(y_i(t; x + z_*\delta))$ is uniformly bounded. Let τ_9 be the first time $y(t; x)$ hits set $\{\|y\| \leq C_9\}$. We claim that, uniformly in $t \in [0, \tau_9]$,

$$(46) \quad \limsup_{\delta \downarrow 0} \frac{1}{\delta} [\xi_i(t; x + z_*\delta, z) - \xi_i(t; x, z)] \leq \epsilon C_*,$$

where $C_* > 0$ is a universal constant. Indeed, let $t_1 \in S(\tau_9; x)$ be the first (smallest) switching point of trajectory $y(\cdot; x)$. To be concrete, let us assume $t_1 > 0$. (The case t_1 is treated analogously.) For a given δ , we define a *switching interval* $[\theta_1^*, \theta_1^{**}]$ associated with t_1 as follows: θ_1^* is the minimum of t_1 and those switching points of $y(\cdot; x + z\delta)$ that are within distance $\epsilon\delta$ from t_1 ; similarly, θ_1^{**} is the maximum of t_1 and those switching points of $y(\cdot; x + z\delta)$ that are within distance $\epsilon\delta$ from t_1 . Obviously, $\theta_1^{**} - \theta_1^* \leq 2\epsilon\delta$. In the interval $[0, \theta_1^*]$, $\xi(t; x + z_*\delta, z) = \xi(t; x, z)$, because they are governed by the ODE with *same* matrix u^m . Within the switching interval, the ODEs for $\xi(t; x + z_*\delta, z)$ and $\xi(t; x, z)$ may have a different matrix u^m , but there is only a finite number of those matrices; therefore, in $[\theta_1^*, \theta_1^{**}]$, $\|\xi(t; x + z_*\delta, z) - \xi(t; x, z)\|$ can increase at most by $C_{11}\|\xi(\theta_1^*; x, z)\|\epsilon\delta$, where C_{11} is a universal constant. We then consider the second switching point t_2 and the associated switching interval $[\theta_2^*, \theta_2^{**}]$. Note that between the first and second switching intervals, both $\xi(t; x + z_*\delta, z)$ and $\xi(t; x, z)$ are again governed by the ODE with same matrix u^m ; therefore the difference $\xi(t; x + z_*\delta, z) - \xi(t; x, z)$ is governed by the same ODE, and therefore in the interval $[\theta_1^{**}, \theta_2^*]$ the value of $\|\xi(t; x + z_*\delta, z) - \xi(t; x, z)\|$ can increase at most by a factor given by a universal constant $C_{12} > 0$ (by (36)). At the end of the switching interval $[\theta_2^*, \theta_2^{**}]$, the first-order component of $\|\xi(t; x + z_*\delta, z) - \xi(t; x, z)\|$ is upper bounded by

$$C_{12}C_{11}\|\xi(\theta_1^*; x, z)\|\epsilon\delta + C_{11}\|\xi(\theta_2^*; x, z)\|\epsilon\delta.$$

We consider the third switching point, and so on. We see that the first-order component of $\|\xi(t; x + z_*\delta, z) - \xi(t; x, z)\|$ will be upper bounded by $C_*\epsilon\delta$, for a sufficiently large universal C_* . (There will be also higher order terms δ^ℓ , $\ell \geq 2$, with uniformly bounded coefficients.) This proves claim (46).

By Lemma 6(i), $\tau_9 \leq C_{14}\|x\|$, for a universal constant $C_{14} > 0$. (We can always choose $C_9 \geq \alpha$, and then $C_{14} = T$.) Then, $\tau_9 \leq C_{14}(g(x) + C_{15})$, where C_{15} may depend on the parameter $C > 0$ of function g .

Now, for all sufficiently small δ , the integral of the term (45),

$$\begin{aligned} & \int_0^\infty \frac{1}{\delta} f'(y_i(t; x + z_*\delta)) [\xi_i(t; x + z_*\delta, z) - \xi_i(t; x, z)] dt = \\ & \int_0^{\tau_9} \frac{1}{\delta} f'(y_i(t; x + z_*\delta)) [\xi_i(t; x + z_*\delta, z) - \xi_i(t; x, z)] dt, \end{aligned}$$

because $f'(y_i(t; x + z_*\delta)) = 0$ for $t \geq T$. The absolute value of the latter integral is upper bounded by

$$C_{14}(g(x) + C_{15})C_*\epsilon = C_{14}C_*\epsilon g(x) + C_{14}C_{15}C_*\epsilon.$$

The constants C_{14} and C_* are universal, while C_{15} depends on C , which depends on C_9 , which depends on ϵ . It remains to choose ϵ small enough so that $C_{14}C_*\epsilon < C_1$. Then the value of $C_{14}C_{15}C_*\epsilon$, plus the corresponding upper bound on the integral of (44), gives constant C_2 . \square

Proof of Theorem 10, case $\mu_{22} = \mu_{12}$. This case is treated the same way as $\mu_{12} \neq \mu_{22}$, with the following modifications. If there is *no* switching point $t \in S(\tau_9; x)$, associated with equality $y_1(t) + y_2(t) = b$, then the proof is unchanged. Suppose there is a switching point $t \in S(\tau_9; x)$, associated with equality $y_1(t) + y_2(t) = b$. Then, in the notation of the proof of Lemma 6, we must have $t \geq t'$, and by the observation we made in that proof, t is the last switching point, and therefore it is the only switching point associated with equality $y_1(t) + y_2(t) = b$. Moreover, all the properties we established in the $\mu_{22} \neq \mu_{12}$ case proof, still apply to all switching points before t . After time t , the process stays within the domain \mathcal{X}_0 , and therefore $(d/dt)[y_1(t) + y_2(t)] = -\mu_{22}[y_1(t) + y_2(t)]$. In particular, in a small neighborhood of time t , $(d/dt)[y_1(t) + y_2(t)] \leq -(b/2)\mu_{22} < 0$. These facts imply that the switching interval, corresponding to switching time t , is such that its end points are within $C_{40}\delta$ from t , for some universal constant $C_{40} > 0$. This means that the contribution of this last switching interval, as well as of the remaining time interval up to the time τ_9 , into the integral of (45), is upper bounded by a universal constant $C_{41} > 0$. \square

6. Generalization of the N -system. Theorem 2, along with its proof, easily extend to the generalization of N -system, shown in Figure 3, in the Halfin-Whitt regime. The system has two customer types and arbitrary number of server pools. There is exactly one server pool that is flexible, i.e. can serve both types. (On Figure 3, it is the pool in the middle.) Each of the remaining pools is dedicated to service of either type 1 or 2. (The two pools on the left in the figure are dedicated to type 1, while the two pools on the right – to type 2.) Each customer type has absolute preference for its dedicated server pools, in some fixed priority order, over the flexible pool. In the flexible pool, the absolute preemptive priority is given to one of the types.

The key features that the generalized system shares with the N -system are that there are two customer types and only one flexible server pool, which can be shared by the customers of different types. These features are exploited in Section 5.2, where we estimated second derivatives of the Lyapunov function. (We note again that all results in Section 5.1, which concern with first derivatives, hold for far more general systems, e.g. those under LAP discipline [14, 15].) The behavior of the DFLs for the generalized system is more complicated, simply because the number of state space domains can be very large. However, as in the N -system, after a finite time all dedicated server pools stay fully occupied, which means that the DFL dynamics depends only on “what happens” in the flexible pool. Consequently, our analysis goes through with very minor adjustments.

7. Discussion. In this paper we address the problem of tightness of stationary distributions, and the limit interchange, for flexible multi-pool service systems in the Halfin-Whitt regime. The behavior of such systems can be very complicated, which makes the problem challenging. This is, in particular, due to the difficulty of constructing Lyapunov functions. We propose an approach which uses a single Lyapunov function, defined as an integral functional of the drift-based fluid limits (DFL) $y(\cdot): G(x) = \int_0^\infty g(y(t))dt$, $y(0) = x$. The problem then reduces to studying the (first and second) derivatives of a DFL – and the corresponding integral $G(x)$ – on the initial state x . We apply this approach to show the tightness property for the N -model under a priority discipline.

Both the approach and many parts of our analysis are quite generic and might be applicable to other models as well. In this respect, note that there is a lot of flexibility in choosing the “distance”

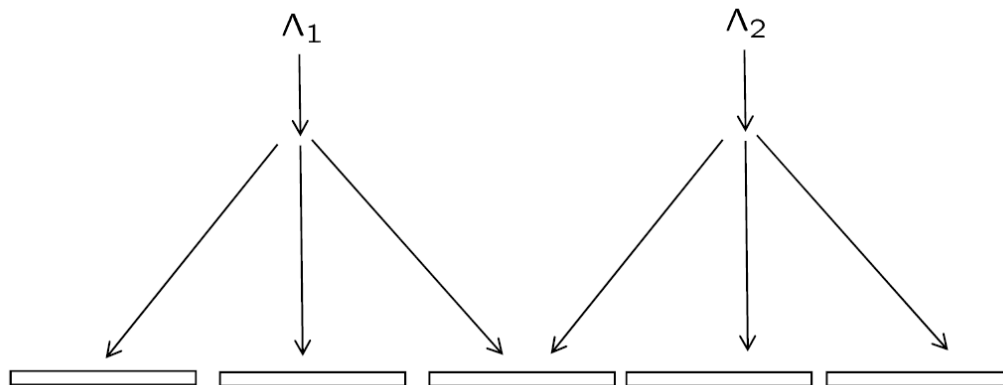


FIG 3. A more general system.

function $g(\cdot)$. It might also be possible to combine the approach with other approaches. For example, a Lyapunov function of the type we consider could be defined and applied on a subspace, if it could be shown by other means that the stationary distributions concentrate (in appropriate sense) on that subspace. Exploring these directions may be a subject of future research.

REFERENCES

- [1] AKSIN, Z., ARMONY, M., AND MEHROTRA, V. (2007). The modern call-center: A multi-disciplinary perspective on operations management research. *Production and Operations Management, Special Issue on Service Operations in honor of John Buzacott (ed. G. Shanthikumar and D. Yao)* **16**, 6, 655–688.
- [2] BUDHIRAJA, A. AND LEE, C. (2009). Stationary distribution convergence for generalized Jackson networks in heavy traffic. *Math. Oper. Res.* **34**, 45–56.
- [3] DAI, G. J., DIEKER, A. B., AND GAO, X. (2014). Validity of heavy-traffic steady-state approximations in many-server queues with abandonment. *Preprint at <http://arxiv.org/abs/1306.5346>*.
- [4] DIEKER, A. B. AND GAO, X. (2012). Positive recurrence of piecewise ornstein-uhlenbeck processes and common quadratic lyapunov functions. *Annals of Applied Probability*.
- [5] GAMARNIK, D. AND GOLDBERG, D. (2013). Steady-state GI/GI/n queue in the Halfin-Whitt regime. *Annals of Applied Probability*.
- [6] GAMARNIK, D. AND MOMCILOVIC, P. (2008). Steady-state analysis of a multiserver queue in the halfin-whitt regime. *Advances in Applied Probability* **40**, 548–577.
- [7] GAMARNIK, D. AND STOLYAR, A. L. (2012). Multiclass multiserver queueing system in the halfin-whitt heavy traffic regime. asymptotics of the stationary distribution. *Queueing Systems* **71**, 25–51.
- [8] GAMARNIK, D. AND ZEEVI, A. (2006). Validity of heavy traffic steady-state approximations in generalized jackson networks. *The Annals of Applied Probability* **16**, 56–90.
- [9] GANS, N., KOOLE, G., AND MANDELBAUM, A. (2003). Telephone call centers: Tutorial, review, and research prospects. *Manufacturing & Service Operations Management* **5**, 79–141.
- [10] GURVICH, I. (2013). Validity of heavy-traffic steady-state approximations in multiclass queueing networks: The case of queue-ratio disciplines. *Mathematics of Operations Research*.
- [11] GURVICH, I. AND WHITT, W. (May 2009). Queue-and-idleness-ratio controls in many-server service systems. *Mathematics of OR* **34**, 2, 363–396.
- [12] HALFIN, S. AND WHITT, W. (1981). Heavy-traffic limits for queues with many exponential servers. *Operations Research* **29**, 3, 567–588.
- [13] SCHONLEIN, M. AND WIRTH, F. (2012). On converse lyapunov theorems for fluid network models. *Queueing Systems* **70**, 339–367.
- [14] STOLYAR, A. L. (2013). Diffusion scale tightness of invariant distributions of a large-scale flexible service system. [arXiv:1301.5838](https://arxiv.org/abs/1301.5838).

- [15] STOLYAR, A. L. AND YUDOVINA, E. (2012). Tightness of invariant distributions of a large-scale flexible service system under a priority discipline. *Stochastic Systems* **2**, 2, 381–408.
- [16] STOLYAR, A. L. AND YUDOVINA, E. (2013). Systems with large flexible server pools: Instability of “natural” load balancing. *Annals of Applied Probability* **23**, 5, 2099–2138.
- [17] YE, H. Q. AND CHEN, H. (2001). Lyapunov method for the stability of fluid networks. *Operations Research Letters* **28**, 125–136.

MURRAY HILL, NJ
E-MAIL: stolyar@research.bell-labs.com